

# Bayesian sparse multiple regression for simultaneous rank reduction and variable selection

Antik Chakraborty

Department of Statistics, Texas A&M University, College Station  
3143 TAMU, TX 77843-3143, USA  
antik@stat.tamu.edu

Anirban Bhattacharya

Department of Statistics, Texas A&M University, College Station  
3143 TAMU, TX 77843-3143, USA  
anirbanb@stat.tamu.edu

Bani K. Mallick

Department of Statistics, Texas A&M University, College Station  
3143 TAMU, TX 77843-3143, USA  
bmallick@stat.tamu.edu

## Abstract

We develop a Bayesian methodology aimed at estimating low rank and row sparse matrices in a high dimensional multivariate response linear regression model. Starting with a full rank matrix and thus avoiding any prior specification on the rank, we let our estimate shrink towards the space of low rank matrices using continuous shrinkage priors. For selecting rows we propose a one step post processing scheme derived from putting group lasso penalties on the rows of the coefficient matrix with default choice of tuning parameters. We then provide an adaptive posterior estimate of the rank using a novel optimization function achieving dimension reduction in the covariate space. We exhibit the performance of the proposed methodology in an extensive simulation study and a real data example.

**Key Words:** Bayesian; High dimension; Low rank; Reduced rank; Dimension reduction; Variable selection.

**Short title:** Bayesian sparse multi-task learner

# 1 Introduction

Studying the relationship between multiple response variables and a set of predictors has broad applications ranging from bioinformatics, econometrics, time series analysis to growth curve models. The least squares solution in a linear multiple response regression problem is equivalent to performing separate least squares on each of the responses (Anderson, 1984) and thus ignores any potential dependence among the responses. Moreover, it lacks interpretability due to the large number of parameters involved in the model. Dimension reduction techniques such as principal component analysis, canonical correlation, factor models have gained wide popularity in multivariate exploratory data analysis due mainly to their ease of interpretation. In the context of multiple response regression, a popular technique to achieve parsimony and interpretability is known as reduced rank regression, wherein a low rank structure of the coefficient matrix is considered. The study of such procedures date back to Anderson (1951), Izenman (1975); a thorough review of the developments in the last century is documented in the monograph by Velu & Reinsel (2013). Although many results exist about the asymptotic properties of reduced rank estimators, formal statistical determination of the rank remains difficult with even fixed number of covariates and large sample size. The problem is substantially harder when the number of covariates is large and thus limits the application of classical reduced rank procedures to routine modern statistical applications such as genomics, financial time series etc.(Bunea et al., 2011). This motivated a series of developments aimed at simultaneously estimating the rank and handle high dimensional covariates. One unifying theme of this line of work relies on penalizing matrices with high ranks which is achieved either by penalizing the singular values of the coefficient matrix (Chen et al., 2013; Yuan et al., 2007) or by penalizing the rank itself (Bunea et al., 2011). Theoretical evaluations of these estimators showed that they adapt to the oracle convergence rate when the true coefficient matrix is of low rank (Bunea et al., 2011). At the same time several authors noted convergence rates can be improved further if one accounts for selecting variables (Tibshirani, 1996; Yuan & Lin, 2006) when the coefficient matrix has many zero rows and low rank. Methods that simultaneously handle low rank and row sparse matrices were then proposed (Bunea et al., 2012; Chen & Huang, 2012; Yuan et al., 2007) penalizing matrices with high rank and dense rows. Uncertainty characterization in the estimation of parameters obtained from these procedures is currently not available

to the best of our knowledge.

The first fully systematic Bayesian treatment of reduced rank regression was carried out in Geweke (1996), where conditioned on the knowledge of the rank, independent Gaussian shrinkage priors were assumed on the elements of the coefficient matrix. Geweke (1996) suggested computing Bayes factors for determining the rank which are typically not mathematically tractable and computationally burdensome in high dimensions. Simplified analytic approximations of the Bayes factor, such as the Savage-Dickey ratio (Dickey, 1971), exist in the literature; Marin & Robert (2010) pointed out a measure theoretic issue with such procedures and proposed a modification. The problem of choosing the rank is not unique to reduced rank regression and is ubiquitous in situations involving low rank decomposition with factor models being a prominent example. Lopes & West (2004) placed a prior on the number of factors and proposed a reversible jump algorithm (Green, 1995) for posterior computation; see also Godsill (2001) for related methods. However, reversible jump MCMC is usually expensive in terms of computation time and mixing. Bhattacharya & Dunson (2011) instead proposed to increasingly shrink the factors starting with a conservative upper bound and adaptively collapse redundant columns inside their MCMC algorithm. Recent advancements in Bayesian reduced rank regression have taken a similar approach in starting with a full rank coefficient matrix and increasingly shrink along the column index induced by stochastically decreasing variances (Alquier, 2013; Babacan et al., 2011; Lim & Teh, 2007; Salakhutdinov & Mnih, 2008). However, it is not clear how choices of hyper parameters involved in these priors may influence the estimation of the rank.

From a Bayesian point of view, a natural way to select variables in a single-response regression framework is to use point mass mixture priors (George & McCulloch, 1993; Scott et al., 2010) which allow a subset of the regression coefficients to be exactly zero. These priors were also used in multiple response regression problems by several authors (Bhadra & Mallick, 2013; Brown et al., 1998; Lucas et al., 2006; Wang, 2010) where independent point mass mixture priors were placed on the rows of the coefficient matrix. Various Gibbs sampling strategies to implement these priors are discussed in Dellaportas et al. (2000). Posterior inference with such priors involve a stochastic search over the exponentially growing model space and is computationally expensive even in moderate dimensions. On the other hand, by applying continuous shrinkage priors (Polson & Scott, 2010) one can mimic the operating characteristics of point mass mixtures while avoiding the computational

bottleneck. The scale mixture of Gaussian representation of these priors renders straightforward Gibbs update of the regression coefficients in a block, with priors in Bhattacharya et al. (2015) and Carvalho et al. (2010) having proven rapid mixing and convergence (Khare & Hobert, 2013; Pal & Khare, 2014). However, such block updates from a high dimensional Gaussian distribution can be very expensive (Bhattacharya et al., 2016). Also, these priors do not achieve variable selection. Hahn & Carvalho (2015) proposed posterior summary aimed at selecting variables with high predictive power to remedy this deficiency. Such posterior summaries have also been proposed in Bondell & Reich (2012) and Kundu et al. (2013).

In this article we simultaneously address the problems of dimension reduction and variable selection in reduced rank models. To the best of our knowledge, no existing Bayesian methodology handles these two aspects at the same time. We develop our methodology with three main building blocks. First, instead of specifying a prior distribution on the rank, we shrink the entire coefficient matrix towards the space of low-rank and row-sparse matrices using continuous shrinkage priors, resulting in sufficient computational gains relative to priors that aim to select variables and/or the rank. We develop an exact sampling scheme to update the entire coefficient matrix in a block leading to cheaper computational cost and better mixing. Posterior summaries of the coefficient matrix, for example the pointwise mean or median, are then processed via two independent screening procedures: one which achieves row-sparsity, another which reduces dimensionality via row sparsity. One of the key features of our post-processing schemes is to come up with careful choices of tuning parameters based on the posterior, resulting in a procedure which is completely free of tuning parameters. We call the resulting procedure Bayesian sparse multi-task learner (BSML). Matlab implementation of the proposed method can be found online at <https://github.com/antik015/Simultaneous-Dimension-Reduction-and-Variable-Selection>.

## 2 Methodology

### 2.1 Model and Prior Specification

Suppose, for each observational unit  $i = 1, \dots, n$ , we have a multivariate response  $y_i \in \mathbb{R}^q$  on  $q$  variables of interest, along with information on  $p$  possible predictors  $x_i \in \mathbb{R}^p$ , a subset of which are assumed to be important in predicting the  $q$  responses. Let  $X \in \mathbb{R}^{n \times p}$  denote the design matrix whose  $i$ th row is  $x_i^T$ , and  $Y \in \mathbb{R}^{n \times q}$  with the  $i$ th row as  $y_i^T$ . The multivariate linear regression

model is,

$$Y = XC + E, \quad E = (e_1^T, \dots, e_n^T)^T, \quad (1)$$

where we follow standard practice to center the response and exclude the intercept term. Our main motivation is the case where  $q < n \ll p$ , although the method trivially applies to  $p < n$  settings as well. We assume that the rows of the error matrix are independent with  $e_i \sim N(0, \Sigma)$ . Depending on the application,  $\Sigma$  can assume different structures; we shall restrict attention to the case  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_q^2)$  in this article, noting that extensions to auto-regressive or low rank structures can be accommodated easily.

In reduced rank regression problems with  $\text{rank}(C) = r \leq \min(p, q)$ ,  $C$  admits a rank decomposition  $C = B_* A_*^T$ , where  $B_* \in \mathbb{R}^{p \times r}$  and  $A_* \in \mathbb{R}^{q \times r}$ . While it is possible to treat  $r$  as a parameter and assign it a prior distribution inside the hierarchical formulation, posterior inference on  $r$  either requires calculation of marginal likelihoods, which is intractable for the present problem, or resorting to complicated reversible jump Markov chain Monte Carlo algorithms. As an alternative, we write  $C = BA^T$ , where  $B = (B_* \mid O^{(q-r) \times q}) \in \mathbb{R}^{p \times q}$  and  $A = (A_* \mid O^{q \times (q-r)}) \in \mathbb{R}^{q \times q}$ , by appending zero columns to  $A_*$  and  $B_*$ . We work in a parameter expanded framework (Liu & Wu, 1999) and aim to approximate  $A$  and  $B$  by full rank matrices by shrinking out the redundant columns. Restricting  $A$  to be in the Stiefel manifold of orthogonal matrices forces  $B$  to have the same row sparsity as the coefficient matrix (Chen & Huang, 2012). Uniform priors over the Stiefel manifold may then be used and computational strategies for posterior inference with such priors were developed in Hoff (2007). We found in our numerical experiments that posterior updates involving uniform priors resulted in relatively slow computation time. To this end, we assign independent Gaussian priors on the entries of  $A$ , and assign shrinkage priors on the columns of  $B$  to encourage column sparsity. Our hierarchical prior is,

$$Y = XBA^T + E, \quad e_i \sim N(0, \Sigma), \quad (2)$$

$$b_{jh} \mid \lambda_{jh}, \tau_h \sim N(0, \lambda_{jh}^2 \tau_h^2), \quad \lambda_{jh} \sim f, \quad \tau_h \sim g, \quad (3)$$

$$a_{hk} \sim N(0, 1), \quad (4)$$

$$\pi(\sigma_h^2) \propto \sigma_h^{-2}, \quad h = 1, \dots, q, \quad (5)$$

for  $j = 1, \dots, p$ ,  $h = 1, \dots, q$ , where  $f$  and  $g$  are densities on the positive real line. Choices of  $f$  and  $g$  crucially impact the performance in sparse situations (Pati et al., 2014; Polson & Scott, 2010). We set  $f$  and  $g$  to be the Half-Cauchy distribution which leads to the Horseshoe prior (Carvalho et al., 2010). In a normal means problem, the posterior with a horseshoe prior adapts to the unknown sparsity and concentrates around the true parameter value at optimal minimax rate (van der Pas et al., 2014). While extensions of more complicated models are lacking in general barring a few exceptions (Pati et al., 2014), empirical studies (Bhattacharya et al., 2016) suggest that similar results can be expected.

## 2.2 Posterior Computation

Exploiting the conditional conjugacy of the proposed prior, we develop a straightforward and efficient Gibb's sampler to update the model parameters in (2) from their full conditional distributions. We use vectorization to update parameters in blocks. Specifically, in what follows, we will make multiple usage of the following identity. For matrices  $\Phi_1, \Phi_2, \Phi_3$  with appropriate dimensions, and  $\text{vec}(A)$  denoting column-wise vectorization, we have,

$$\text{vec}(\Phi_1 \Phi_2 \Phi_3) = (\Phi_3^T \otimes \Phi_1) \text{vec}(\Phi_2) = (\Phi_3^T \Phi_2^T \otimes I_k) \text{vec}(\Phi_1) \quad (6)$$

where the matrix  $\Phi_1$  has  $k$  rows and  $\otimes$  denotes the Kronecker product.

The Gibbs sampler cycles through the following step, sampling parameters from their full conditional distributions:

**Step 1.** To sample  $B \mid -$ , use (6) to vectorize  $Y = XBA^T + E$  to obtain,

$$y = (X \otimes A)\beta + e, \quad (7)$$

where  $\beta = \text{vec}(B^T) \in \mathbb{R}^{pq \times 1}$ ,  $y = \text{vec}(Y^T) \in \mathbb{R}^{nq \times 1}$ , and  $e = \text{vec}(E^T) \sim N_{nq}(0, \tilde{\Sigma})$  with  $\tilde{\Sigma} = \text{diag}(\Sigma, \dots, \Sigma)$ . Multiplying both sides of (7) by  $\tilde{\Sigma}^{-1/2}$  yields  $\tilde{y} = \tilde{X}\beta + \tilde{e}$  where  $\tilde{y} = \tilde{\Sigma}^{-1/2}y$ ,  $\tilde{X} = \tilde{\Sigma}^{-1/2}(X \otimes A)$  and  $\tilde{e} = \tilde{\Sigma}^{-1/2}e \sim N_{nq}(0, I_{nq})$ . Then the full conditional distribution of  $\beta \mid -$  is given by  $\beta \mid - \sim N_{pq}(\Omega_B^{-1} \tilde{X}^T \tilde{y}, \Omega_B^{-1})$ , where  $\Omega_B = (\tilde{X}^T \tilde{X} + \Lambda^{-1})$  with  $\Lambda = \text{diag}(\lambda_{11}^2 \tau_1^2, \dots, \lambda_{1q}^2 \tau_q^2, \dots, \lambda_{p1}^2 \tau_1^2, \dots, \lambda_{pq}^2 \tau_q^2)$ .

Naively sampling from the full conditional of  $\beta$  has complexity  $O(p^3 q^3)$  which becomes highly expensive for moderate values of  $p$  and  $q$ . Bhattacharya et al. (2016) recently developed an algo-

rithm to sample from a class of structured multivariate normal distributions as in the full conditional of  $\beta$  whose complexity scales linearly in the ambient dimension. We adapt the algorithm in Bhattacharya et al. (2016) as follows:

- (i) Sample  $u \sim \mathcal{N}(0, \Lambda)$  and  $\delta \sim \mathcal{N}(0, \mathbf{I}_{nq})$  independently.
- (ii) Set  $v = \tilde{X}u + \delta$ .
- (iii) Solve  $(\tilde{X}\Lambda\tilde{X}^\top + \mathbf{I}_{nq})w = (\tilde{y} - v)$  to obtain  $w$ .
- (iv) Set  $\beta = u + \Lambda\tilde{X}^\top w$ .

The salient feature is that one only requires matrix multiplications and linear system solvers to implement the above algorithm, and no matrix decomposition is required. It follows from standard results on complexity (Golub & van Loan, 1996) of matrix operations that the above steps have a combined complexity of  $O(q^3 \max\{n^2, p\})$ , a substantial improvement over  $O(p^3 q^3)$  when  $p \gg n$ .

**Step 2.** To sample  $A \mid -$ , once again vectorize  $Y = XBA^\top + E$ , but this time we use the equality of the first and the third terms in (6) to obtain,

$$y = (XB \otimes \mathbf{I}_q)a + e, \quad (8)$$

where  $e$  and  $y$  are the same as in step 1, and  $a = \text{vec}(A) \in \mathbb{R}^{q^2 \times 1}$ . The full conditional posterior distribution  $a \mid - \sim \mathcal{N}(\Omega_A^{-1}X_*^\top \tilde{y}, \Omega_A^{-1})$ , where  $\Omega_A = (X_*^\top X_* + \mathbf{I}_{q^2})$ ,  $X_* = \tilde{\Sigma}^{-1/2}(XB \otimes \mathbf{I}_{q^2})$  and  $\tilde{y} = \tilde{\Sigma}^{-1/2}y$ . To sample from the full conditional of  $a$ , we use the following steps from Rue (2001). Compute the Cholesky decomposition  $(X_*^\top X_* + \mathbf{I}_{q^2}) = LL^\top$ . Solve the following system of equations:  $Lv = X_*^\top \tilde{y}$ ,  $L^\top m = v$ , and  $L^\top w = z$ , where  $z \sim \mathcal{N}(0, \mathbf{I}_{q^2})$ . Finally obtain a sample as  $a = m + w$ .

**Step 3.** To sample  $\sigma_h^2 \mid -$ , observe that  $\sigma_h^2 \mid - \sim \text{Inverse-Gamma}(n/2, S_h/2)$  independently across  $h$ , where  $S_h = \{Y_h - (XBA^\top)_h\}^\top \{Y_h - (XBA^\top)_h\}$ , with  $\Phi_h$  denoting the  $h$ th column of a matrix  $\Phi$ .

**Step 4.** The global and local scale parameters  $\lambda_{jh}$ 's and  $\tau_h$ 's have independent conditional posteriors across  $j$  and  $h$ , which can be sampled via a slice sampling scheme provided in the online supplement to Polson et al. (2014). We illustrate the sampling technique for a generic local shrinkage parameter  $\lambda_{jh}$ ; a similar scheme works for  $\tau_h$ . Setting  $\eta_{jh} = \lambda_{jh}^{-2}$ , the slice sampler proceeds by sampling  $u_{jh} \mid \eta_{jh} \sim \text{Unif}(0, 1/(1 + \eta_{jh}))$  and then sampling  $\eta_{jh} \mid u_{jh} \sim \text{Exp}(2\tau_h^2/b_{jh}^2)\mathbf{I}\{\eta_{jh} < (1 - u_{jh})/u_{jh}\}$ , a truncated exponential distribution.

### 2.3 Post processing for variable selection

Point estimates  $\hat{C}_B$ , posterior mean or element-wise posterior median, of  $C$  obtained from the Gibb's sampler above can be readily used for prediction purposes, with a natural quantification of predictive uncertainty. However, the continuous nature of our prior implies that such point estimates will be non-sparse and full rank with probability one. Relying on the fact that  $\hat{C}_B$  will increasingly concentrate around the true parameter (Alquier, 2013) under a wide range of shrinkage priors, we develop a two step automated procedure to arrive at an estimate of  $C$  which is row sparse and rank reduced. The first step is to zero out rows of  $\hat{C}_B$  aimed at removing irrelevant predictors. To achieve such row sparsity, we consider an utility function based approach which minimizes prediction error in future observations with an added penalty for model complexity. Specifically, using the posterior mean  $\bar{C}$  as our choice for  $\hat{C}_B$ , the row sparse estimate  $\hat{C}_R$  for  $C$  is given by

$$\hat{C}_R = \arg \min_{\Gamma} \| X\bar{C} - X\Gamma \|_F^2 + \sum_{j=1}^p \mu_j \| \Gamma^{(j)} \|_2 \quad (9)$$

where  $\Phi^{(j)}$  represents the  $j^{th}$  row of a matrix  $\Phi$ , and the  $\mu_j$ 's are predictor specific regularization parameters. The group lasso (Yuan & Lin, 2006) penalty  $\sum_{j=1}^p \mu_j \| \Gamma^{(j)} \|_2$  implies row sparsity in the resulting solution. The main difference from a usual group lasso estimator is achieved by replacing the noisy data matrix  $Y$  by the fit  $X\bar{C}$  from the shrinkage procedure. For a derivation of the objective function in (9) from a utility function perspective (Hahn & Carvalho, 2015), refer to the Appendix.

Hahn & Carvalho (2015) suggested inspection of a type of posterior summary plots to arrive at sparse estimates in the univariate regression context. It is not immediately clear how to extend the summary plots technique to matrix-valued coefficients. Instead, we provide a soft thresholding scheme with default choices of thresholds. Specifically, we set the sub-gradient of (9) with respect to  $\Gamma^{(j)}$  to zero and replace  $\| \Gamma^{(j)} \|$  by a data dependent quantity to obtain the soft thresholding estimate,

$$\hat{C}_R^{(j)} = \frac{1}{X_j^T X_j} \left( 1 - \frac{\mu_j}{2 \| X_j^T R_j \|} \right)_+ X_j^T R_j, \quad (10)$$

where  $R_j$  is the residual matrix obtained after regressing  $X\bar{C}$  on  $X$  leaving out the  $j^{th}$  predictor,  $R_j = X\bar{C} - \sum_{k \neq j} X_k \hat{C}_R^{(k)}$ . We use  $\bar{C}$  as our initial estimate and make a pass through each variable and update the initial estimate according to (10). Derivation of (10) is deferred to the Appendix.



While the  $p$  tuning parameters  $\mu_j$  can be chosen by cross-validation, the computation cost explodes with  $p$  to search over a grid in  $p$  dimensions. We instead recommend the default choice  $\mu_j = \|\overline{C}^{(j)}\|^{-2}$  which has good empirical operating characteristics. We stress here that with this default choice, one iteration suffices and no iterative algorithm is required to obtain  $\widehat{C}_R$ .

## 2.4 Post processing for rank estimation

To obtain an estimate of the rank, we threshold the singular values of  $X\widehat{C}_R$ , with  $\widehat{C}_R$  obtained from (10). In situations where row sparsity is not warranted,  $\overline{C}$  can be used instead of  $\widehat{C}_R$ . Letting  $s_1, \dots, s_q$  denote the singular values of  $X\widehat{C}_R$ , consider

$$\arg \min_{\nu_1, \dots, \nu_q} \sum_{h=1}^q (s_h - \nu_h)^2 + \omega \sum_{h=1}^q \mathbf{I}(\nu_h \neq 0). \quad (11)$$

The solution of (11) can be explicitly written as a hard thresholding scheme:  $\nu_h = s_h \mathbf{I}(s_h > \omega)$ . We estimate the rank as the number of nonzero thresholded singular values, that is,  $\hat{r} = \sum_{h=1}^q \mathbf{I}(s_h > \omega)$ . It is important for prediction purpose that we operate on the prediction space  $\{XC \in \mathbb{R}^{n \times q} : \text{rank}(C) \leq k\}$  instead of the estimation space  $\{C \in \mathbb{R}^{p \times q} : \text{rank} \leq k\}$ . A row sparse and rank reduced estimate of  $C$ , denoted  $\widehat{C}_{RR}$ , is obtained as the best rank  $\hat{r}$  approximation to  $\widehat{C}_N$ , where  $\widehat{C}_N$  is the matrix component with non-zero rows of  $\widehat{C}_R$ :  $\widehat{C}_R = (\widehat{C}_N \mid O)$ . This is done using the singular value decomposition of  $\widehat{C}_N$  truncated to the first  $\hat{r}$  terms and then appending the zero rows obtained in the previous step. The estimator  $\widehat{C}_{RR}$  so derived is called the Bayesian sparse multitask learning estimator.

Letting  $s_0$  be the largest singular value of  $Y - X\widehat{C}_R$ , we use  $s_0$  as a default choice for the threshold parameter  $\omega$ . Intuitively, the idea is to retain the singular values of  $\widehat{C}_R$  which are above the noise level, with  $Y - X\widehat{C}_R$  a natural estimator of the noise matrix, whose largest singular value is then used to quantify the noise level.

## 3 Simulation Results

We performed a thorough simulation study to assess the performance of the proposed method across different settings. For all our simulation settings the sample size  $n$  was fixed at 100. We chose 4 different  $(p, q)$  combinations:  $(300, 10)$ ,  $(500, 10)$ ,  $(200, 30)$ ,  $(1000, 12)$ . For each of these combinations the set of true predictors was  $J = \{1, 2, \dots, 10\}$ , with  $p_0 = |J| = 10$ . The data

were generated from the model  $Y = XB_*A_*^T + E$ , with  $B_* \in \mathbb{R}^{p \times r}$  and  $A_* \in \mathbb{R}^{r \times r}$ , with the rank  $r \in \{3, 5, 7\}$ . The elements of the first  $p_0$  rows of  $B_*$  and all entries of  $A_*$  were independently generated from  $N(0, 1)$ , while the last  $(p - p_0)$  rows of  $B_*$  were set to 0. We let  $C_0 = B_*A_*^T$  denote the true coefficient matrix.

The columns of the design matrix  $X$  were independently generated from  $N(0, \Sigma_X)$ . We considered two cases,  $\Sigma_X = I_p$ , and  $\Sigma_X = (\sigma_{ij}^X)$ ,  $\sigma_{jj}^X = 1$ ,  $\sigma_{ij}^X = 0.5$  for  $i \neq j$ . Each row of the matrix  $E$  was generated from a multivariate normal distribution with diagonal covariance matrix having diagonal entries uniformly chosen between 0.5 and 1.75.

Letting  $A^{(t)}$  and  $B^{(t)}$  denote the value of  $A$  and  $B$  at the  $t^{th}$  iteration of the MCMC chain we define  $C^{(t)} = B^{(t)}A^{(t)T}$  and the posterior mean and the posterior median are respectively the pointwise mean and median of  $T$  such iterations leaving out the initial burnin samples. We compare the performance of our row and rank selected estimator in terms of the Mean Square Error (MSE) and Mean Square Prediction Error (MSPE) with the following definitions: for any estimator  $\hat{C}$  of  $C_0$ ,

$$\text{MSE} = \|\hat{C} - C_0\|_F^2 / pq \quad \text{and} \quad \text{MSPE} = \|X\hat{C} - XC_0\|_F^2 / nq. \quad (12)$$

We considered two choices of  $\hat{C}$ :  $\hat{C}_{RR}$  obtained using the posterior mean and the pointwise posterior median. We denote the MSE of the posterior mean by MSE and that of the posterior median by  $\text{MSE}_{\text{me}}$  and similar notations are used for the MSPE. For relative comparison we considered the method due to Yuan et al. (2007), hereafter referred as Sparse Partial Least Squares (SPLS).

Table 1 reports the average MSE ( $\times 10^{-4}$ ) and MSPE across 50 replicates. The first row for each rank contains the average value of the estimated rank for the two methods. The SPLS estimator of the rank is the one for which the minimum cross validation error is achieved. In our simulations we used the default 10 fold cross validation in the `cv.spls` function from the R package `spls`. The results of our proposed rank estimator are highly accurate, whereas corresponding numbers for SPLS show it achieves moderate dimension reduction. Moreover, SPLS estimates of the rank are highly sensitive to the number of response variable involved in the problem. The estimates increase significantly from column 2, 4, 6, 8 to 10 in table 1. To that end, BSML is very adaptive and does not show any tendency of overfitting at the cost of increased prediction error. Our rank estimator by design assumes a low rank coefficient matrix and is expected to perform better in cases where

$(r, q) = (3, 10), (3, 12)$  compared to the cases where  $(r, q) = (7, 10), (7, 12)$ . In spite of that our numerical results indicate that, this estimator works reasonably well even when the true coefficient matrix is moderately dense in terms of rank, namely the case where  $(r, q) = (7, 10), (7, 12)$ .

For all the settings considered,  $\hat{C}_{RR}$  obtained from the posterior mean and the posterior median clearly outperforms SPLS in terms of estimation and prediction error. BSML estimates adapt to the underlying sparsity and dimension excellently. For some cases the BSML improves estimation error over SPLS by more than 100 times. One can make similar observations about the prediction error also. For example see columns 13 and 14 of table 1.

We now discuss the variable selection performance of the proposed method. We compare BSML and SPLS in terms of their sensitivity and specificity. The sensitivity is defined as the ratio of the correctly selected variables and the true number of variables and the specificity is defined as the ratio between the number of variables removed and the true number of noise variables. Table 2 summarizes the variable selection performance of  $\hat{C}_{RR}$  under the same set of simulation setups as in Table 1. As in table 1 we consider the same choices of  $\hat{C}_{RR}$ .

Considering the sensitivity and specificity results in table 2 there is little difference in BSML and SPLS at least when the dimension is moderate. For example, SPLS shows high sensitivity for  $p = 300, 200, 500$  but the sensitivity decreases significantly for  $p = 1000$ . But for BSML we see that the sensitivity results do not change with increasing number of covariates. Methods which have high specificity tend to over-select at the cost of increased false positives. BSML maintains high specificity while correctly removing the noise variables which is not the case for SPLS.

One more remark we want to make here is that the good performance of SPLS in variable selection comes at the cost of very moderate dimension reduction. We considered the case  $(n, p, q, r) = (100, 300, 10, 3)$  to compare the variable selection performance of BSML and SPLS with ranks fixed at 3, 6, 9 with 20 replicates. Every time BSML selected the correct set variables but for SPLS the average percentage of false positives were 10 and it failed to select 5 signal variables on average with rank fixed at the true value 3. With rank 6 the false positive percentage was 5 and SPLS missed 2 true variables. As such, with SPLS there is apparently a trade off between dimension reduction and variable selection. However, that is not the case with BSML; for BSML rank misspecification does not result in poor variable selection.

Table 1: Estimation and predictive performance of the proposed method (BSML) versus SPLS across different simulation settings. We report the average estimated rank ( $\hat{r}$ ), Mean Square Error ( $\times 10^{-4}$ ) and Mean Square Predictive Error for  $\hat{C}_{RR}$  obtained using the posterior mean (MSE, MSPE) and posterior median (MSE<sub>me</sub>, MSPE<sub>me</sub>) across 50 replications. For each setting the true number of signals were 10 and sample size was 100. For each combination of  $(p, q, r)$  the columns of the design matrix were generated from  $N(0, \Sigma_X)$ . Two different choices of  $\Sigma_X$  was considered.  $\Sigma_X = I_p$  (independent) and  $\Sigma_X = (\sigma_{ij}^X), \sigma_{jj}^X = 1, \sigma_{ij}^X = 0.5$  for  $i \neq j$  (correlated). Better estimates of rank and minimum MSE and MSPE achieved for each setting is in bold characters.

Rank	Measures	(p,q)															
		(300,10)				(500,10)				(200,30)				(1000,12)			
		Independent		Correlated		Independent		Correlated		Independent		Correlated		Independent		Correlated	
		BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS
3	$\hat{r}$	<b>3.0</b>	8.8	<b>2.9</b>	10	<b>3.0</b>	9.7	<b>3.0</b>	8.8	<b>3.0</b>	7.9	<b>3.0</b>	9.4	<b>3.2</b>	9.4	<b>3.4</b>	8.9
	MSE	<b>4</b>	7	<b>6</b>	13	<b>3</b>	7	<b>5</b>	30	<b>3</b>	14	<b>5</b>	15	<b>3</b>	50	<b>3</b>	38
	MSE <sub>me</sub>	3		5		2		4		3		5		1		2	
	MSPE	<b>0.11</b>	0.20	<b>0.11</b>	0.23	<b>0.17</b>	0.22	<b>0.15</b>	0.34	<b>0.07</b>	0.25	<b>0.06</b>	0.21	<b>0.35</b>	4.19	<b>0.30</b>	1.51
	MSPE <sub>me</sub>	0.08		0.09		0.09		0.11		0.05		0.05		0.12		0.16	
5	$\hat{r}$	<b>4.9</b>	9.9	<b>4.8</b>	9.5	<b>4.9</b>	9.9	<b>4.8</b>	9.8	<b>5</b>	9.7	<b>4.9</b>	12.2	<b>5.1</b>	9.9	<b>5.1</b>	9.9
	MSE	<b>4</b>	10	<b>8</b>	40	<b>3</b>	10	<b>6</b>	24	<b>5</b>	690	<b>9</b>	617	<b>2</b>	108	<b>4</b>	129
	MSE <sub>me</sub>	4		7		2		5		5		8		1		3	
	MSPE	<b>0.13</b>	0.3	<b>0.14</b>	0.53	<b>0.17</b>	0.41	<b>0.20</b>	0.38	<b>0.11</b>	3.8	<b>0.09</b>	4.6	<b>0.32</b>	9.54	<b>0.32</b>	4.63
	MSPE <sub>me</sub>	0.10		0.12		0.11		0.15		0.09		0.08		0.14		0.19	
7	$\hat{r}$	<b>6.8</b>	10	<b>6.7</b>	10	<b>6.8</b>	10	<b>6.7</b>	9.7	<b>6.9</b>	10.3	<b>6.9</b>	15.8	<b>6.8</b>	10.2	<b>6.6</b>	11.5
	MSE	<b>4</b>	10	<b>8</b>	90	<b>3</b>	20	<b>5</b>	49	<b>6</b>	1162	<b>10</b>	1128	<b>2</b>	195	<b>4</b>	261
	MSE <sub>me</sub>	4		7		2		5		5		9		1		3	
	MSPE	<b>0.14</b>	0.43	<b>0.13</b>	0.86	<b>0.16</b>	0.72	<b>0.16</b>	0.92	<b>0.12</b>	10.81	<b>0.11</b>	9.01	<b>0.32</b>	16.70	<b>0.31</b>	7.44
	MSPE <sub>me</sub>	0.11		0.12		0.11		0.13		0.10		0.10		0.15		0.18	

Table 2: Variable selection performance after the proposed post processing step is reported across the same simulation settings as in 1. The sensitivity and specificity reported is the average across 50 replications. For each  $r = 3, 5, 7$  the third row reports the standard deviation ( $\times 10^{-3}$ ) of the sensitivity and specificity across 50 replicates. Selection accuracy is compared with SPLS. As in table 1 for each case maximum sensitivity and specificity is in bold characters.

Rank	Measures	(p,q)															
		(300,10)				(500,10)				(200,30)				(1000,12)			
		Independent		Correlated		Independent		Correlated		Independent		Correlated		Independent		Correlated	
		BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS	BSML	SPLS
3	Sens.	<b>1</b>	<b>1</b>	<b>1</b>	0.98	<b>1</b>	<b>1</b>	<b>1</b>	0.98	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	0.59	<b>1</b>	0.73
	Spec.	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	<b>0.98</b>	<b>0.98</b>	<b>0.99</b>	<b>0.99</b>	0.97	<b>0.99</b>	0.98	<b>1</b>	0.96	<b>0.99</b>
	Std.	(0, 4)	(0, 8)	(0, 8)	(1, 11)	(0, 4)	(0, 5)	(0, 20)	(7, 20)	(0, 5)	(0, 10)	(0, 20)	(0, 5)	(0, 20)	(90, 0)	(0, 30)	(99, 10)
5	Sens.	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	0.57	<b>1</b>	0.72	<b>1</b>	0.56	<b>1</b>	0.68
	Spec.	<b>0.99</b>	0.96	<b>0.99</b>	0.95	<b>0.99</b>	<b>0.99</b>	0.98	<b>0.99</b>	0.99	<b>1</b>	0.95	<b>0.96</b>	0.99	<b>1</b>	0.98	<b>0.99</b>
	Std.	(0, 7)	(0, 4)	(0, 5)	(0, 6)	(0, 7)	(0, 9)	(0, 20)	(0, 16)	(0, 7)	(10, 0)	(0, 20)	(86, 50)	(0, 5)	(69, 4)	(0, 10)	(92, 8)
7	Sens.	<b>1</b>	<b>1</b>	<b>1</b>	0.98	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	0.51	<b>1</b>	0.68	<b>1</b>	0.54	<b>1</b>	0.74
	Spec.	<b>0.99</b>	<b>0.99</b>	<b>0.99</b>	0.97	<b>0.99</b>	0.97	<b>0.98</b>	<b>0.98</b>	0.99	<b>1</b>	0.95	<b>0.99</b>	0.99	<b>1</b>	<b>0.98</b>	<b>0.98</b>
	Std.	(0, 2)	(0, 7)	(0, 5)	(3, 12)	(0, 7)	(0, 8)	(0, 20)	(0, 18)	(0, 10)	(30, 4)	(0, 30)	(98, 49)	(0, 50)	(80, 0)	(0, 20)	(70, 42)

## 4 Yeast Cell Cycle Data

Identifying transcription factors (TF) which are responsible for cell cycle regulation is an important scientific problem (Chun & Keleş, 2010). The yeast cell cycle data from Spellman et al. (1998) contains information from three different experiments on mRNA levels of 800 genes on an  $\alpha$ -factor based experiment. The response variable is the amount of transcription (mRNA) which was measured every 7 minutes in a period of 119 minutes, a total of 18 measurements (Y) covering two cell cycle periods. The ChIP-chip data from Lee et al. (2002) on chromatin immunoprecipitation contains the binding information of the 800 genes for 106 TFs (X). We analyze this data obtained from the R package `spls` which has the above information completed for 542 genes.

When applied to this data, our method selected 33 of the 106 TF's compared to 48 for multi-variate SPLS and 69 for SRRR Chen & Huang (2012). We ran 20 independent runs of the Gibbs sampler with different starting points to check the stability of our variable selection procedure. Each time the same set of 33 variables were selected. Out of the 21 scientifically verified TFs (Wang et al., 2007) our method selected 14 TFs. For SPLS and SRRR these numbers are 14 and 16 respectively. Although SRRR selects more number correct TFs, SRRR also selected 53 irrelevant TFs.

The effect of the 21 confirmed TF's as estimated from our procedure are drawn in figure 1. 10 additional TF's that regulate cell cycle were identified by Lee et al. (2002), out of which 3 TF's were selected by our proposed method. We do not enumerate common TF's that are selected by SPLS, SRRR and the BSML for the sake of brevity. Some examples include, ACE2, SWI4, SWI5, SWI6 etc. A direct comparison of the estimated effects of these TF's from the three methods mentioned earlier reveals the similar periodic patterns obtained from the three methods. However, it is to be noted that in particular for SWI4, SRRR estimate effects are much smaller in magnitude than what is estimated from our methodology. This is not surprising given the presence of local scale parameters in our model which allow for different shrinkage for each entry of the coefficient matrix. Whereas, for SRRR specially, each row is shrunk towards zero by the same amount. Periodicity in the estimated effects is not at all surprising since the mRNA levels were measured for 2 hours covering two cell cycles.

Our automatic rank detection technique estimated a rank of 1 which is significantly different

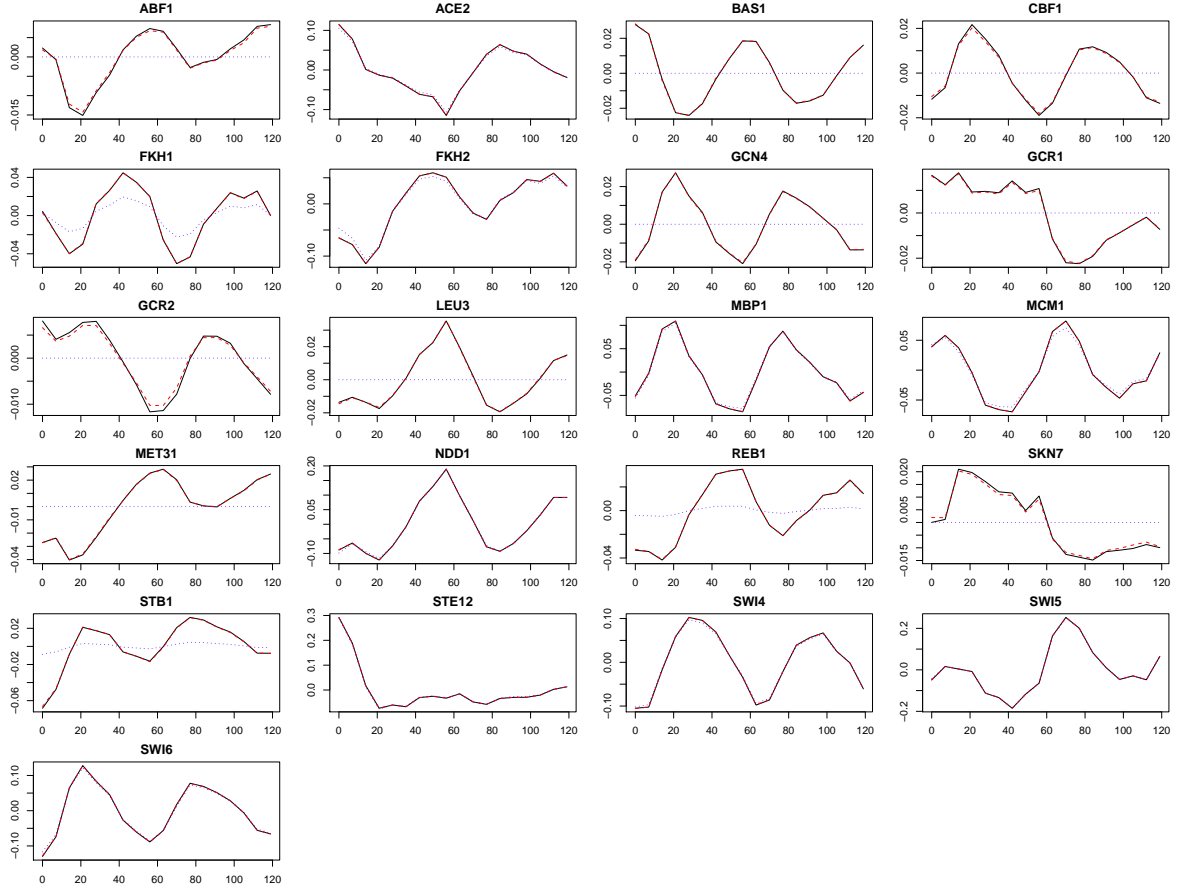
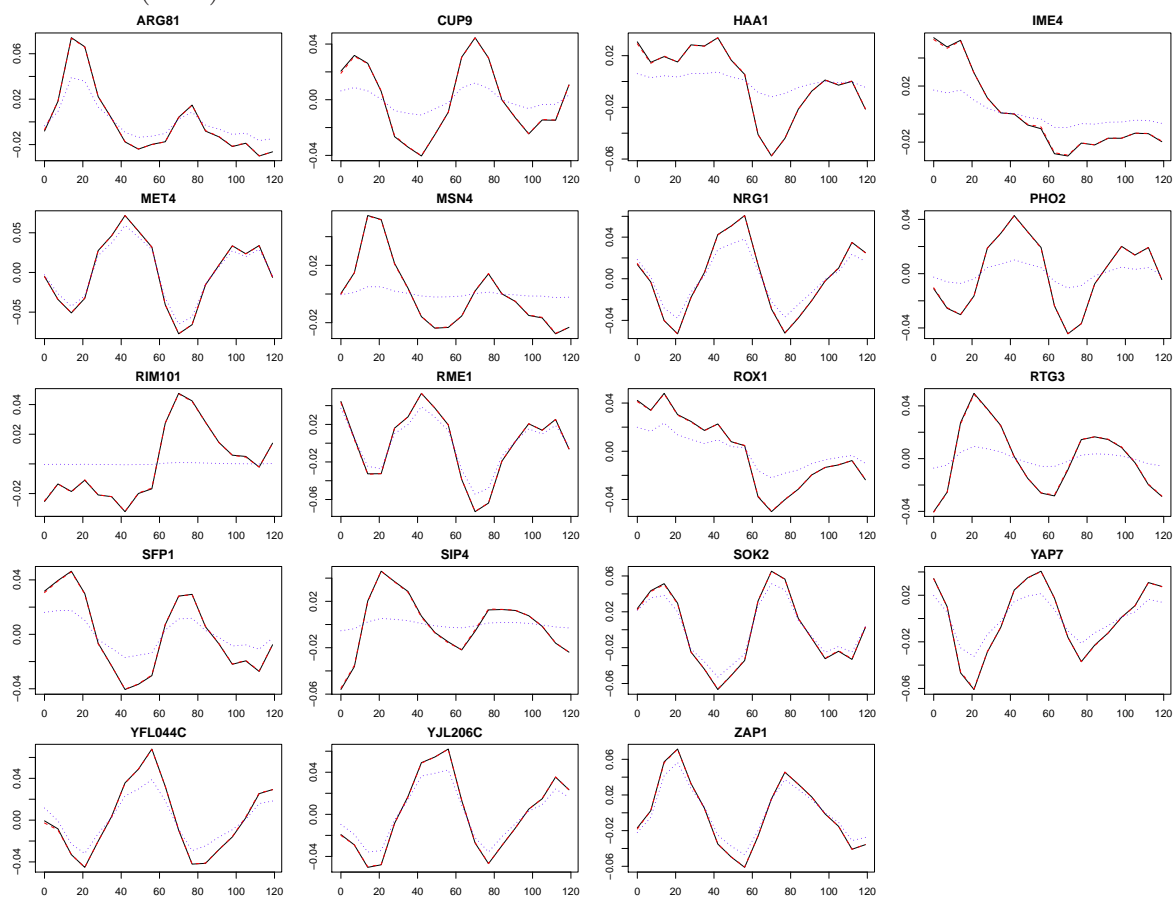


Figure 1: Estimated effects of the 21 scientifically verified TF's. On  $y$ -axis are the estimated coefficients and on  $x$ -axis are the measuring time in minutes. The estimated coefficients are plotted for the posterior mean (solid black line), posterior median (red dashed line) and sparse estimate after post processing (blue dotted line). 14 of the verified TF's were selected using the proposed method.

Figure 2: Estimated effects of the TF's selected by the proposed method that are not experimentally verified. The lines represent the same as in figure 1. Of these SIP4 and YAP7 was also selected by Yuan & Lin (2006).



from SRRR (4) and SPLS (8). The proposed rank estimator depends on the validity of the linear model assumed between the responses and the covariates and also the error covariance structure. We suspect that although the linear model assumption might be valid here, the error covariance matrix  $\Sigma$  should not be assumed identity since the mRNA measurements were taken over a period of 120 minutes with equal intervals. However, the singular values of  $Y - XC_R$  for the yeast data showed a significant drop in magnitude after the first four values which agrees with the findings in Chen & Huang (2012).

## A Appendix

### A.1 Derivation of (9)

Suppose  $Y^* \in \mathbb{R}^{n \times q}$  be  $n$  future observations with design points  $X$  so that given  $C$ ,  $Y^*$  can be decomposed into  $Y^* = XC + E^*$  where  $E^*$  has the same distribution as  $E$  in (1). We define the utility function in terms of loss of predicting these  $n$  new future observations. To encourage sparsity in rows of a coefficient matrix  $\Gamma$  that balances the prediction we add a group lasso penalty (Yuan & Lin, 2006) to this utility function. We define the utility function as,

$$\mathcal{L}(Y^*, \Gamma) = \|Y^* - X\Gamma\|_F^2 + \sum_{j=1}^p \mu_j \|\Gamma^{(j)}\|_2 \quad (13)$$

where the  $p$  tuning parameters  $\{\mu_j\}_{j=1}^p$  control the penalty for selecting each predictor variable and  $\Phi^{(j)}$  represents the  $j^{th}$  row of any matrix  $\Phi$ . Intuitively we want  $\mu_j$  to be small if the  $j^{th}$  predictor is important and vice versa. The expected risk,  $\mathbb{E}\{\mathcal{L}(Y^*, \Gamma)\}$ , after integrating over the space of all such future observations given  $C$  and  $\Sigma$ , is

$$\mathcal{L}(\Gamma, C, \Sigma) = q \operatorname{tr}(\Sigma) + \|XC - X\Gamma\|_F^2 + \sum_{j=1}^p \mu_j \|\Gamma^{(j)}\|_2. \quad (14)$$

Finally we take expectation of this quantity with respect to  $\pi(C, \Sigma \mid Y, X)$  and drop the constant terms to obtain (9).

### A.2 Derivation of (10)

We let  $\Phi_j$  and  $\Phi^{(j)}$  denote the  $j^{th}$  column and row of a generic matrix  $\Phi$  and unless otherwise mentioned  $\|x\|$  denotes the Euclidean norm of a vector  $x$ . Using the subgradient of (9) with



respect to  $\Gamma^{(j)}$  (Friedman et al., 2007), we have

$$2X_j^T(X\Gamma - X\overline{C}) + \mu_j\alpha_j = 0, \quad j = 1, \dots, p, \quad (15)$$

where  $\alpha_j = \Gamma^{(j)} / \|\Gamma^{(j)}\|$  if  $\|\Gamma^{(j)}\| \neq 0$  and  $\alpha_j = 0$  when  $\|\Gamma^{(j)}\| = 0$ . For  $\Gamma^{(j)} = 0$  we can rewrite (15) as,  $2X_j^T(\sum_{k \neq j} X_k \Gamma^{(k)} - X\overline{C}) + \mu_j\alpha_j = 0$  which imply that  $\alpha_j = -\frac{2}{\mu_j}X_j^T R_j$ , where  $R_j$  is the residual matrix obtained after regressing  $X\overline{C}$  on  $X$  leaving out the  $j^{th}$  predictor,  $R_j = X\overline{C} - \sum_{k \neq j} X_k \Gamma^{(k)}$ . We can use this to set  $\Gamma^{(j)}$  to zero: if  $\alpha_j < 1$  set  $\Gamma^{(j)} = 0$ . Otherwise we have  $2X_j^T(X_j\Gamma^{(j)} - R_j) + \mu_j\frac{\Gamma^{(j)}}{\|\Gamma^{(j)}\|} = 0$ . Solving for  $\Gamma^{(j)}$  in the above equation we then get,

$$\Gamma^{(j)} = \left( X_j^T X_j + \frac{\mu_j}{2\|\Gamma^{(j)}\|} \right)^{-1} X_j^T R_j. \quad (16)$$

This solution is dependent on the unknown quantity  $\|\Gamma^{(j)}\|$ . However, taking norm on both sides in (16) we get a value of  $\|\Gamma^{(j)}\|$  which does not involve any unknown quantities:  $\|\Gamma^{(j)}\| = (\|X_j^T R_j\| - \mu_j/2) / X_j^T X_j$ . Substituting this in (16) we get,  $\Gamma^{(j)} = \frac{1}{X_j^T X_j} \left( 1 - \frac{\mu_j}{2\|X_j^T R_j\|} \right) X_j^T R_j$ .

Finally, combining the case when  $\Gamma^{(j)} = 0$ , we have (10).

## References

- ALQUIER, P. (2013). Bayesian methods for low-rank matrix estimation: Short survey and theoretical study. In *Algorithmic Learning Theory*. Springer.
- ANDERSON, T. (1984). Multivariate statistical analysis. *Wiley and Sons, New York, NY*.
- ANDERSON, T. W. (1951). Estimating linear restrictions on regression coefficients for multivariate normal distributions. *The Annals of Mathematical Statistics*, 327–351.
- BABACAN, S. D., MOLINA, R. & KATSAGGELOS, A. K. (2011). Variational bayesian super resolution. *Image Processing, IEEE Transactions on* **20**, 984–999.
- BHADRA, A. & MALLICK, B. K. (2013). Joint high-dimensional bayesian variable and covariance selection with an application to eqtl analysis. *Biometrics* **69**, 447–457.
- BHATTACHARYA, A., CHAKRABORTY, A. & MALLICK, B. K. (2016). Fast sampling with gaussian scale mixture priors in high-dimensional regression. *Biometrika*.

- BHATTACHARYA, A. & DUNSON, D. (2011). Sparse Bayesian infinite factor models. *Biometrika* **98**, 291–306.
- BHATTACHARYA, A., PATI, D., PILLAI, N. S. & DUNSON, D. B. (2015). Dirichlet–laplace priors for optimal shrinkage. *Journal of the American Statistical Association* **110**, 1479–1490.
- BONDELL, H. D. & REICH, B. J. (2012). Consistent high-dimensional bayesian variable selection via penalized credible regions. *Journal of the American Statistical Association* **107**, 1610–1624.
- BROWN, P. J., VANNUCCI, M. & FEARN, T. (1998). Multivariate bayesian variable selection and prediction. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **60**, 627–641.
- BUNEA, F., SHE, Y. & WEGKAMP, M. H. (2011). Optimal selection of reduced rank estimators of high-dimensional matrices. *The Annals of Statistics* , 1282–1309.
- BUNEA, F., SHE, Y., WEGKAMP, M. H. et al. (2012). Joint variable and rank selection for parsimonious estimation of high-dimensional matrices. *The Annals of Statistics* **40**, 2359–2388.
- CARVALHO, C., POLSON, N. & SCOTT, J. (2010). The horseshoe estimator for sparse signals. *Biometrika* **97**, 465–480.
- CHEN, K., DONG, H. & CHAN, K.-S. (2013). Reduced rank regression via adaptive nuclear norm penalization. *Biometrika* , ast036.
- CHEN, L. & HUANG, J. Z. (2012). Sparse reduced-rank regression for simultaneous dimension reduction and variable selection. *Journal of the American Statistical Association* **107**, 1533–1545.
- CHUN, H. & KELEŞ, S. (2010). Sparse partial least squares regression for simultaneous dimension reduction and variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72**, 3–25.
- DELLAPORTAS, P., FORSTER, J. J. & NTZOUFRAS, I. (2000). Bayesian variable selection using the gibbs sampler. *BIostatistics-BASEL-* **5**, 273–286.

- DICKEY, J. M. (1971). The weighted likelihood ratio, linear hypotheses on normal location parameters. *The Annals of Mathematical Statistics* , 204–223.
- FRIEDMAN, J., HASTIE, T., HÖFLING, H., TIBSHIRANI, R. et al. (2007). Pathwise coordinate optimization. *The Annals of Applied Statistics* **1**, 302–332.
- GEORGE, E. I. & MCCULLOCH, R. E. (1993). Variable selection via gibbs sampling. *Journal of the American Statistical Association* **88**, 881–889.
- GEWEKE, J. (1996). Bayesian reduced rank regression in econometrics. *Journal of econometrics* **75**, 121–146.
- GODSILL, S. J. (2001). On the relationship between markov chain monte carlo methods for model uncertainty. *Journal of Computational and Graphical Statistics* **10**, 230–248.
- GOLUB, G. H. & VAN LOAN, C. F. (1996). *Matrix Computations*. John Hopkins University Press, 3rd ed.
- GREEN, P. J. (1995). Reversible jump markov chain monte carlo computation and bayesian model determination. *Biometrika* **82**, 711–732.
- HAHN, P. R. & CARVALHO, C. M. (2015). Decoupling shrinkage and selection in bayesian linear models: a posterior summary perspective. *Journal of the American Statistical Association* **110**, 435–448.
- HOFF, P. D. (2007). Gibbs sampling of the matrix bingham-von mises-fisher distribution, with an application to protein interaction networks .
- IZENMAN, A. J. (1975). Reduced-rank regression for the multivariate linear model. *Journal of multivariate analysis* **5**, 248–264.
- KHARE, K. & HOBERT, J. P. (2013). Geometric ergodicity of the bayesian lasso. *Electronic Journal of Statistics* **7**, 2150–2163.
- KUNDU, S., BALADANDAYUTHAPANI, V. & MALLICK, B. K. (2013). Bayes regularized graphical model estimation in high dimensions. *arXiv preprint arXiv:1308.3915* .

- LEE, T. I., RINALDI, N. J., ROBERT, F., ODOM, D. T., BAR-JOSEPH, Z., GERBER, G. K., HANNETT, N. M., HARBISON, C. T., THOMPSON, C. M., SIMON, I. et al. (2002). Transcriptional regulatory networks in *saccharomyces cerevisiae*. *science* **298**, 799–804.
- LIM, Y. J. & TEH, Y. W. (2007). Variational bayesian approach to movie rating prediction. In *Proceedings of KDD cup and workshop*, vol. 7. Citeseer.
- LIU, J. S. & WU, Y. N. (1999). Parameter expansion for data augmentation. *Journal of the American Statistical Association* **94**, 1264–1274.
- LOPES, H. F. & WEST, M. (2004). Bayesian model assessment in factor analysis. *Statistica Sinica*, 41–67.
- LUCAS, J., CARVALHO, C., WANG, Q., BILD, A., NEVINS, J. & WEST, M. (2006). Sparse statistical modelling in gene expression genomics. *Bayesian Inference for Gene Expression and Proteomics* **1**, 0–1.
- MARIN, J.-M. & ROBERT, C. P. (2010). On resolving the savage–dickey paradox. *Electronic Journal of Statistics* **4**, 643–654.
- PAL, S. & KHARE, K. (2014). Geometric ergodicity for bayesian shrinkage models. *Electronic Journal of Statistics* **8**, 604–645.
- PATI, D., BHATTACHARYA, A., PILLAI, N. S., DUNSON, D. et al. (2014). Posterior contraction in sparse bayesian factor models for massive covariance matrices. *The Annals of Statistics* **42**, 1102–1130.
- POLSON, N. G. & SCOTT, J. G. (2010). Shrink globally, act locally: sparse bayesian regularization and prediction. *Bayesian Statistics* **9**, 501–538.
- POLSON, N. G., SCOTT, J. G. & WINDLE, J. (2014). The bayesian bridge. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76**, 713–733.
- RUE, H. (2001). Fast sampling of gaussian markov random fields. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **63**, 325–338.

- SALAKHUTDINOV, R. & MNIH, A. (2008). Bayesian probabilistic matrix factorization using markov chain monte carlo. In *Proceedings of the 25th international conference on Machine learning*. ACM.
- SCOTT, J. G., BERGER, J. O. et al. (2010). Bayes and empirical-bayes multiplicity adjustment in the variable-selection problem. *The Annals of Statistics* **38**, 2587–2619.
- SPELLMAN, P. T., SHERLOCK, G., ZHANG, M. Q., IYER, V. R., ANDERS, K., EISEN, M. B., BROWN, P. O., BOTSTEIN, D. & FUTCHER, B. (1998). Comprehensive identification of cell cycle-regulated genes of the yeast *saccharomyces cerevisiae* by microarray hybridization. *Molecular biology of the cell* **9**, 3273–3297.
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* , 267–288.
- VAN DER PAS, S., KLEIJN, B., VAN DER VAART, A. et al. (2014). The horseshoe estimator: Posterior concentration around nearly black vectors. *Electronic Journal of Statistics* **8**, 2585–2618.
- VELU, R. & REINSEL, G. C. (2013). *Multivariate reduced-rank regression: theory and applications*, vol. 136. Springer Science & Business Media.
- WANG, H. (2010). Sparse seemingly unrelated regression modelling: Applications in finance and econometrics. *Computational Statistics & Data Analysis* **54**, 2866–2877.
- WANG, L., CHEN, G. & LI, H. (2007). Group scad regression analysis for microarray time course gene expression data. *Bioinformatics* **23**, 1486–1494.
- YUAN, M., EKICI, A., LU, Z. & MONTEIRO, R. (2007). Dimension reduction and coefficient estimation in multivariate linear regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **69**, 329–346.
- YUAN, M. & LIN, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**, 49–67.